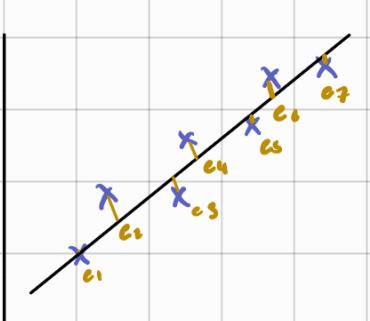


Note : Cornell University (put papers on public domain)
 Learn Prompt engineering course

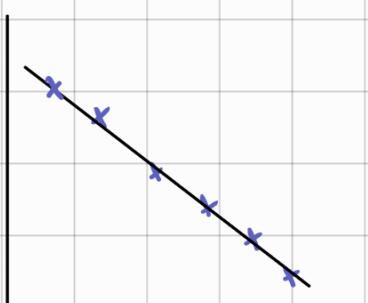
Chat GPT (Generative Pre-trained Transformer Model)

is a LLM that generates text.

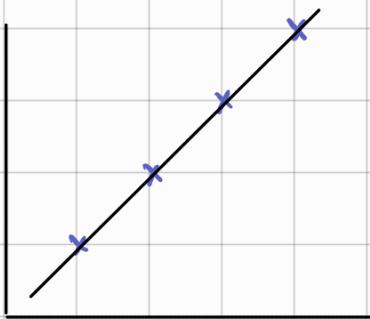
Chat GPT is an optimization model that modifies the answers slightly.



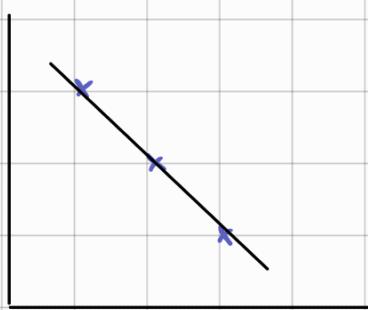
equation of the
model $\hat{y} = \hat{\beta}x + \hat{\alpha}$
 $\min(\sum_{i=1}^n e_i^2)$



$S_{ny} \approx -0.8 < 0$



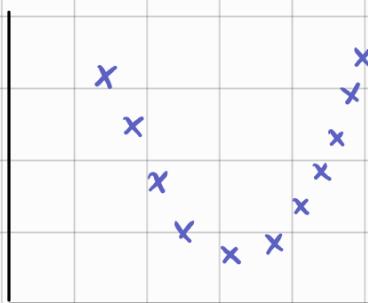
$S_{ny} = 1 > 0$



$S_{ny} = -1 < 0$



$S_{ny} = 0$



$S_{ny} = 0$ (non linear)

$$S_{ny} = \frac{Cov(x, y)}{\sqrt{\text{Var } x \text{ Var } y}}$$

$$= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Linear Regression assumes that the relation between the

variables is linear $y = \beta n + \alpha + \varepsilon \rightarrow \text{error}$

Real Life Events Comprise of

- 1) Uncertainty
- 2) Variability

$$y_i = \alpha + \beta n_i + \varepsilon_i$$

↓ ↓ ↓
intercept slope error
independent variable

error distribution always follows normal dist

Instead of minimizing ε_i , we minimize $\sum \varepsilon_i^2$

Sum of Squared Error: $SS = \min \left(\sum_{i=1}^n (y_i - \beta n_i - \alpha)^2 \right)$

$$\frac{\partial SS}{\partial \alpha} \stackrel{\text{set}}{=} 0 \quad \frac{\partial SS}{\partial \beta} \stackrel{\text{set}}{=} 0$$

see we get optimal: $\hat{\alpha}$ and $\hat{\beta}$
 $y_i = \hat{\beta} n_i + \hat{\alpha}$

For now, we assume: $\hat{\beta} = \frac{S_{xy}}{S_{xx}} = \frac{\text{Cov}(x, y)}{\text{Var}(x)} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

x_i	\bar{x}	y_i	\bar{y}	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
1	3	1	2	-2	-1	2
2		1		-1	-1	1
3		2		0	0	0
4		2		1	0	0
5		4		2	2	4

$$\hat{\beta} = \frac{2+1+4}{4+1+1+4} = \frac{7}{10}$$

$$\hat{\alpha} = 2 - \frac{7}{10} \cdot 3 = \frac{-1}{10}$$

$y_i = 0.7n_i - 0.1$ is the best fitted line!

- * Regression is a tool for investigating the relationship between a dependent variable and one or more indep. variables. It has applications in economics, AI, health, engineering, etc. Before investigating, we study the correlations and plot the scatter plot to study the extent of the linear association.

Simple regression is used to study and model the relationship between one single regressor X and one response Y .

X : Regressor or control, not RV, has uncertainty & variability (regressing y on n).

Y : Dependent random variable

Equation of simple linear regression: $y = \beta_1 n + \beta_0 + \varepsilon$ where ε is the random error, β_0 intercept & β_1 slope.

Assumption: Both X and Y are continuous i.e they cannot be $\{0, 1\}$ but can be $[0, 1]$

- * Physical interpretation of β_0 and β_1 : $\hat{Y} = \beta_0 + \beta_1 X$

we think of \hat{Y} as the profit or price and its unit is y , x 's unit is n so the unit of β_0 is y and the unit of β_1 is y/n . We usually think of X as the investment.

- * Given the data $(x_i, y_i)_{i \in \mathbb{N}}$, we check the scatter plot to find whether the linear model is appropriate or not (if not change the model).

Simple linear regression model assumptions :

- 1) Homoscedasticity Assumption :

ε_i is a random variable with mean 0 and std = σ unknown i.e $\varepsilon_i \sim N(0, \sigma^2)$, $\mathbb{E}(\varepsilon_i) = 0$ $\text{Var}(\varepsilon_i) = \sigma^2$

- 2) $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \Rightarrow \varepsilon_i$ and ε_j are independent

- 3) $\varepsilon_i \stackrel{\text{indep}}{\sim} N(0, \sigma^2)$

stochastic variable non-stochastic variable

We deduce that : $\mathbb{E}(Y_i) = \mathbb{E}(\beta_0 + \beta_1 x_i + \varepsilon_i) = \beta_0 + \beta_1 x_i$

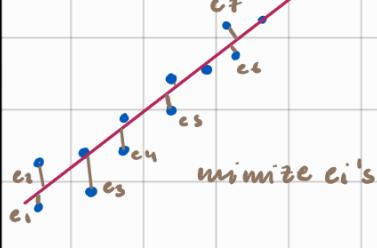
$\text{Var}(Y_i) = \text{Var}(\beta_0 + \beta_1 x_i + \varepsilon_i) = \text{Var}(\varepsilon_i) = \sigma^2$

$Y_i \stackrel{\text{indep}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2)$

Then, if Y_i is not normal, the model fails.

- * Best Fitted β_0 and β_1 :

The line fitted by the least square is the one that makes the sum of squares of the vertical discrepancies as small as possible. Our aim is to estimate β_0 and β_1 so that the sum of squares of all the differences between the observations y_i and the line (fitted) is minimum.



Sum of squared errors (SS) :

$$S = \min \left(\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \right) = \min \left(\sum_{i=1}^n e_i^{n2} \right)$$

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad \varepsilon_i \text{ is estimated by } e_i \text{ (residual error)}$$

y_i	\hat{y}_{iD}	\hat{y}_{iz}
-	-	-
-	-	-

The best is $\min \{ \sum (y_i - \hat{y}_{iD})^2, \sum (y_i - \hat{y}_{iz})^2 \}$

Residual Error : $e_i = y_i - \hat{y}_i$
 true response estimated response

* Least square estimator at $\hat{\beta}_0$ and $\hat{\beta}_1$ must satisfy :

Normal equations :

$$\begin{cases} \frac{\partial S}{\partial \beta_0} \Big|_{\hat{\beta}_0, \hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 & (1) \\ \frac{\partial S}{\partial \beta_1} \Big|_{\hat{\beta}_0, \hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 & (2) \end{cases}$$

normal

see $\hat{\beta}_0$ and $\hat{\beta}_1$ are the solutions of the normal equations.

* Solving the equations by Ordinary least square (OLS) :

$$1) \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \Leftrightarrow n \bar{Y} - n \hat{\beta}_0 - n \hat{\beta}_1 \bar{X} = 0 \\ \Leftrightarrow \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

$$2) \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 \Leftrightarrow \sum_{i=1}^n (y_i - (\bar{Y} - \hat{\beta}_1 \bar{X}) - \hat{\beta}_1 x_i) x_i = 0 \\ \Leftrightarrow \sum_{i=1}^n (y_i - \bar{Y}) x_i - \hat{\beta}_1 (x_i - \bar{X}) x_i = 0$$

$$\sum_{i=1}^n (y_i - \bar{Y})(x_i - \bar{x}) =$$

$$\Leftrightarrow \sum_{i=1}^n (y_i - \bar{Y}) x_i = \hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x}) x_i$$

$$\sum_{i=1}^n (y_i - \bar{Y}) x_i - \bar{n} \sum_{i=1}^n (y_i - \bar{Y}) =$$

$$\Leftrightarrow \hat{\beta}_1 = \sum_{i=1}^n (y_i - \bar{Y}) x_i$$

$$\sum_{i=1}^n (y_i - \bar{y}) u_i = n(\bar{y} - \bar{\bar{y}})$$

$$\sum_{i=1}^n (y_i - \bar{y}) u_i \Leftrightarrow \hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\Leftrightarrow \hat{\beta}_1 = \frac{\text{Cov}(x, y)}{\text{Var}(x)} = \frac{s_{xy}}{s_{xx}} \begin{matrix} \rightarrow \text{without } y_n \\ \rightarrow \text{without } x_n \end{matrix}$$

$$\sum_{i=1}^n (y_i - \bar{y}) \bar{x} = \bar{x} \sum_{i=1}^n y_i - \bar{x} \bar{y} n = n \bar{x} \bar{y} - n \bar{x} \bar{y} = 0$$

$$\sum_{i=1}^n (x_i - \bar{x}) \bar{x} = \bar{x} \sum_{i=1}^n x_i - n \bar{x}^2 = n \bar{x}^2 - n \bar{x}^2 = 0$$

So we added a zero term!

* Properties of Least Square Fit :

1) $\sum_{i=1}^n e_i = 0$ sum of residuals in linear model = 0

comes from the 1st normal equation

2) $\sum_{i=1}^n e_i u_i = 0$ sum of product of residual with regressor = 0

comes from the 2nd normal equation

3) $\sum_{i=1}^n e_i \hat{y}_i = 0$ but $\sum_{i=1}^n e_i y_i \neq 0$

* What if we use MLE method?

$$f(\varepsilon_i) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{\varepsilon_i^2}{2\sigma^2}\right)$$

$$L(z_i) = \prod_{i=1}^n f(z_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{z_i^2}{2\sigma^2}\right) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n z_i^2}$$

$$\mathcal{L}(z_i) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2\right)$$

$$\mathcal{L}^* = \log(\mathcal{L}) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial \mathcal{L}^*}{\partial \beta_0} = -\frac{n}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (-1) = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\frac{\partial \mathcal{L}^*}{\partial \beta_1} = -\frac{n}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (-x_i) = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i = 0$$

$$\frac{\partial \mathcal{L}^*}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 = 0$$

$$\Rightarrow \begin{cases} \hat{\beta}_0 \text{ MLE} = \bar{Y} - \hat{\beta}_1 \bar{X} \\ \hat{\beta}_1 \text{ MLE} = S_{XY} / S_{YY} \\ \hat{\sigma}^2 \text{ MLE} = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \text{ different than } \sigma^2 \text{ OLE} \\ \hat{\sigma}^2 = \mathbb{E}\left(\frac{SS_{RES}}{n-2}\right) \end{cases}$$

* Both $\hat{\beta}_1$ and $\hat{\beta}_0$ are linear estimators because :

$$\cdot \quad \hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{Y}) x_i}{\sum_{i=1}^n (x_i - \bar{x}) x_i} = \sum_{i=1}^n c_i y_i, \quad c_i = \frac{(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

$$\cdot \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

* Both estimators are unbiased as : measure of the goodness

$$E(\hat{\beta}_1) = E\left(\frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}\right)$$

$$= \frac{\sum(x_i - \bar{x}) E(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}$$

$$= \frac{\beta_1 \sum(x_i - \bar{x})^2}{\sum(x_i - \bar{x})^2}$$

$$= \beta_1$$

$$E\left(\frac{\sum(x_i - \bar{x})(\beta_0 + \beta_1 x_i + \varepsilon_i - \bar{y})}{\sum(x_i - \bar{x})^2}\right)$$

$$E\left(\frac{\sum(x_i - \bar{x})(\beta_0 + \varepsilon_i - \bar{y}) + \beta_1 \sum(x_i - \bar{x}) x_i}{\sum(x_i - \bar{x})^2}\right)$$

$$\beta_1 E\left(\frac{\sum(x_i - \bar{x})^2}{\sum(x_i - \bar{x})^2}\right)$$

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

$$\bar{y} = \beta_0 + \beta_1 \bar{x} + \bar{\varepsilon}$$

$$(y_i - \bar{y}) = \beta_1(x_i - \bar{x}) + (\varepsilon_i - \bar{\varepsilon})$$

$$E(y_i - \bar{y}) = \beta_1(x_i - \bar{x}) + 0 \text{ as}$$

$$\varepsilon_i \sim N(0, \sigma^2) \text{ so } \bar{\varepsilon} \sim N(0, \sigma^2/n)$$

$$E(\hat{\beta}_0) = E(\bar{y} - \hat{\beta}_1 \bar{x}) = E(\beta_0 + \beta_1 \bar{x} - \hat{\beta}_1 \bar{x})$$

$$= \beta_0 + \beta_1 \bar{x} - \bar{x} E(\hat{\beta}_1)$$

$$= \beta_0 + \beta_1 \bar{x} - \bar{x} \beta_1$$

$$= \beta_0$$

$$* \text{Var}(\hat{\beta}_1) = \text{Var}\left(\frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}\right)$$

$$= \text{Var}\left(\sum_{i=1}^n c_i y_i\right)$$

$$y_i \stackrel{i.i.d.}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2)$$

$$\downarrow \text{so } \text{Cov}(y_i, y_j) = 0$$

$$= \sum_{i=1}^n c_i^2 \text{Var}(y_i)$$

$$= \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\left(\sum_{i=1}^n (x_i - \bar{x})^2\right)^2} \cdot \sigma^2$$

$$= \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$= \frac{\sigma^2}{S_{xx}}$$

$$\text{Var}(\hat{\beta}_0) = \text{Var}(\bar{Y} - \hat{\beta}_1 \bar{X})$$

$$= \text{Var}(\bar{Y}) + \bar{X}^2 \text{Var}(\hat{\beta}_1) - 2 \text{Cov}(\bar{Y}, \hat{\beta}_1)$$

$$= \frac{\sigma^2}{n} + \bar{X}^2 \cdot \frac{\sigma^2}{S_{xx}} - 2 \text{Cov}\left(\frac{1}{n} \sum_{i=1}^n Y_i, \sum_{i=1}^n c_i Y_i\right)$$

$$= \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{S_{xx}} \right) - \frac{2}{n} \sum_{i=1}^n c_i \text{Var}(Y_i)$$

$$= \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{S_{xx}} \right) - \frac{2\sigma^2}{n} \left(\frac{\sum_{i=1}^n (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \right)$$

$$= \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{S_{xx}} \right) - \frac{2\sigma^2}{n} \left(\frac{n\bar{x} - n\bar{X}}{\sum (x_i - \bar{x})^2} \right)$$

$$= \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{S_{xx}} \right)$$

* Esimerkien ott σ^2 :

$$\mathbb{E}\left(\frac{SS_{\text{Res}}}{n-2}\right) = \sigma^2, \quad \frac{SS_{\text{Res}}}{n-2} \text{ is an unbiased estimator of } \sigma^2$$

$$\begin{aligned} \text{Residual: } \sum_{i=1}^n e_i^2 &= \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \\ &= \sum_{i=1}^n (Y_i - \beta_0 - \hat{\beta}_1 x_i)^2 \\ &= \sum_{i=1}^n (Y_i - \bar{Y} + \hat{\beta}_1 \bar{X} - \hat{\beta}_1 x_i)^2 \\ &= \sum_{i=1}^n ((Y_i - \bar{Y})^2 + \hat{\beta}_1^2 (x_i - \bar{x})^2 - 2\hat{\beta}_1 (x_i - \bar{x})(Y_i - \bar{Y})) \\ &= S_{yy} + \hat{\beta}_1^2 S_{xx} - 2\hat{\beta}_1 S_{xy} \\ &= S_{yy} + \hat{\beta}_1^2 S_{xx} - 2\hat{\beta}_1^2 S_{xx} \\ &= S_{yy} - \hat{\beta}_1^2 S_{xx} \end{aligned}$$

$$\begin{aligned}
 \sum (\beta_0 + \beta_1 x_i)^2 - n(\beta_0 + \beta_1 \bar{x}) &= n\beta_0^2 + 2\beta_0\beta_1 n\bar{x} + \beta_1^2 \sum x_i^2 - n\beta_0^2 \\
 &\quad - 2n\beta_0\beta_1 \bar{x} - n\beta_1^2 \bar{x}^2 \\
 \sum (x_i - \bar{x})^2 &= \beta_1^2 (\sum x_i^2 - \bar{x}^2) \\
 \sum x_i^2 - \sum 2x_i \bar{x} + \sum \bar{x}^2 &= \beta_1^2 \sum (x_i - \bar{x})^2 \\
 \sum x_i^2 - 2n\bar{x} + n\bar{x}^2 &= \beta_1^2 S_{xx} \\
 \sum x_i^2 - \bar{x}^2
 \end{aligned}$$

$$\begin{aligned}
 E(SS_{\text{Res}}) &= E(S_{yy} - \hat{\beta}_1^2 S_{xx}) \\
 &= E(S_{yy}) - S_{xx} E(\hat{\beta}_1^2) \\
 &= E\left(\sum_{i=1}^n (y_i - \bar{y})^2\right) - S_{xx} (\text{Var}(\hat{\beta}_1) + E(\hat{\beta}_1)^2) \\
 &= E\left(\sum_{i=1}^n (y_i^2 + \bar{y}^2 - 2y_i \bar{y})\right) - S_{xx} (\sigma^2/S_{xx} + \beta_1^2) \\
 &= E\left(\sum_{i=1}^n y_i^2 + n\bar{y}^2 - 2\bar{y}\sum_{i=1}^n y_i\right) - \sigma^2 - \beta_1^2 S_{xx} \\
 &= \sum_{i=1}^n E(y_i^2) - nE(\bar{y}^2) - \sigma^2 - \beta_1^2 S_{xx} \\
 &= \sum_{i=1}^n (\text{Var}(y_i) + E(y_i)^2) - n(\text{Var}\bar{y} + E(\bar{y})^2) - \sigma^2 - \beta_1^2 S_{xx} \\
 &= n\sigma^2 + \sum_{i=1}^n (\beta_0 + \beta_1 x_i)^2 - \sigma^2 - n(\beta_0 + \beta_1 \bar{x})^2 - \sigma^2 - \beta_1^2 S_{xx} \\
 &= \sigma^2(n-1) + \beta_1^2 \left(\sum_{i=1}^n (x_i - \bar{x})^2\right) - \sigma^2 - \beta_1^2 S_{xx} \\
 &= \sigma^2(n-1) + \beta_1^2 S_{xx} - \sigma^2 - \beta_1^2 S_{xx} \\
 &= \sigma^2(n-2)
 \end{aligned}$$

$$\text{Hence, } E\left(\frac{SS_{\text{Res}}}{n-2}\right) = \sigma^2, \quad MS_{\text{Res}} = \frac{SS_{\text{Res}}}{n-2}$$

* What is the sampling distribution of MS_{Res} ?

$$y_i \stackrel{iid}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2)$$

$$e_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$SS_{\text{Res}} = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y})^2 = \sum_{i=1}^n \left(\frac{y_i - \beta_0 - \beta_1 x_i}{\sigma}\right)^2 \sim \chi_{n-2}^2$$

$$\text{so } \frac{e_i^2}{\sigma^2} \sim \chi_1^2 \text{ and } \sum_{i=1}^n \frac{e_i^2}{\sigma^2} \sim \chi_{n-2}^2 \text{ restrictions: } \sum e_i = 0$$

normal equations $\sum e_i x_i = 0$

$$\frac{SS_{\text{Res}}}{\sigma^2} \sim \chi_{n-2}^2 \Rightarrow \frac{SS_{\text{Res}}/n-2}{\sigma^2/n-2} = \frac{MS_{\text{Res}}}{\sigma^2} (n-2) \sim \chi_{n-2}^2$$

$$MSR \sim \frac{s^2}{n-2} \chi^2_{n-2}$$

* Statistic test at β_1 : Test at slope coefficient

we want to test if \exists a linear relationship between x_i & y_i .

$$H_0: \beta_1 = 0 \quad \text{vs} \quad H_1: \beta_1 \neq 0$$

↓

↓

no linear relationship \exists a linear relationship

$$\hat{\beta}_1 = \beta_0$$

$$\hat{\beta}_1 = \sum_{i=1}^n c_i y_i, \quad y_i \sim N(\beta_0 + \beta_1 x_i, s^2) \quad \hat{\beta}_1 \sim N(\beta_1, \frac{s^2}{S_{xx}})$$

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{s^2/S_{xx}}} \sim N(0, 1), \text{ we use z-test if } s \text{ is known}$$

$$z = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{s^2/S_{xx}}}, \text{ if } |z| > z_{\alpha/2} \text{ we reject } H_0$$

two tailed

$$\text{If } s^2 \text{ is unknown: } t = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{MS_{RES}/S_{xx}}}, \text{ reject if } |t| > t_{\alpha/2, n-2}$$

↓
do not estimate s^2

$$MS_{RES} = \frac{SS_{RES}}{n-2} = \frac{\sum_{i=1}^n e_i^2}{n-2} = \frac{\sum_{i=1}^n (y_i - \hat{y})^2}{n-2}$$

$$\text{For } \beta_0: t^* = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{MS_{RES}} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}}$$

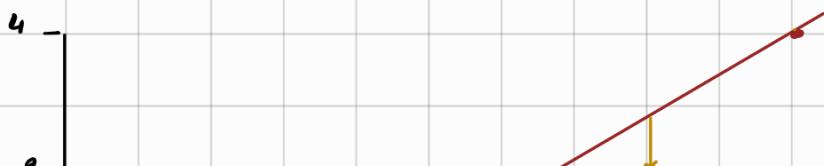
* Example : Show that $\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = 0$ when $\bar{x} = 0$

$$\begin{aligned}
 \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) &= \text{Cov}(\bar{Y} - \hat{\beta}_1 \bar{X}, \sum_{i=1}^n c_i Y_i) \\
 &= \text{Cov}\left(\frac{1}{n} \sum_{i=1}^n Y_i, \frac{\sum (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} Y_i\right) \\
 &= \frac{1}{n} \frac{\sum (x_i - \bar{x})}{S_{xx}} \text{Cov}(Y_i, Y_i) \\
 &= \frac{\text{Var } Y_i}{n S_{xx}} \left(\sum_{i=1}^n x_i - \bar{x} \right) \\
 &= \frac{6^2}{n S_{xx}} (n \bar{x} - n \bar{x}) \\
 &= 0
 \end{aligned}$$

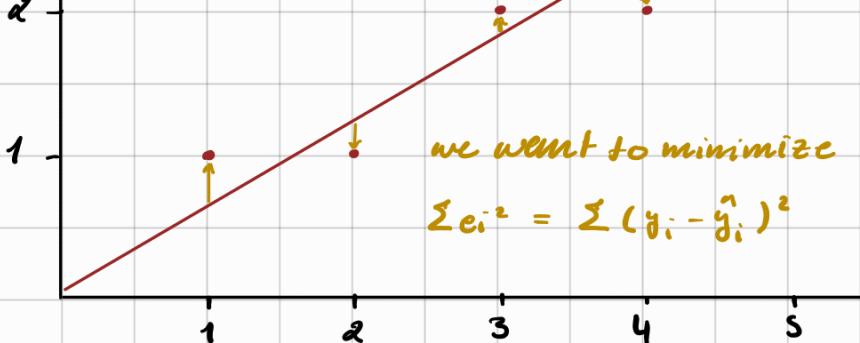
* Example : Suppose you have been hired by a chocolate company as a marketing analyst. They have the following data :

Advertisement cost (10K \$)	Sales (10 ³ units)
1	1
2	1
3	2
4	2
5	4

a) Is there a relationship between sales and advertisement ?



This is a positive



There is a positive association between x_i and y_i as confirmed by the scatter plot and r_{xy}

x_i	\bar{x}	$x_i - \bar{x}$	$(x_i - \bar{x})^2$	s_{xx}	$Var x$	y_i	\bar{y}	$y_i - \bar{y}$	$(y_i - \bar{y})^2$	s_{yy}	$Var y$
1	3	-2	4	10	2	1	2	-1	1	6	$\frac{6}{5}$
2		-1	1			1		-1	1		
3		0	0			2		0	0		
4		1	1			2		0	0		
5		2	4			4		2	4		

$$r_{xy} = \frac{s_{xy}}{\sqrt{s_{xx} s_{yy}}} = \frac{2+1+4}{\sqrt{10} \cdot \sqrt{6}} = \frac{7}{\sqrt{60}} = 0.9$$

$$\hat{\beta}_1 = \frac{s_{xy}}{s_{xx}} = \frac{7}{10} = 0.7 \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 2 - \frac{7}{10} \cdot 3 = -0.1$$

$$\text{Model: } \hat{y} = -0.1 + 0.7x$$

↓ ↓

non sense as companies 0.7 increase in sales per 1\$ (10k \$) at investment sell units even without advertising, we need more data to estimate accurately.

b) Is the relationship significant at 5% level

x_i	y_i	\hat{y}	$e_i = y_i - \hat{y}$	e_i^2	$\sum e_i^2$

$$MS_{RES} = \frac{SS_{RES}}{n-2}$$

1	1	0.4	0.4	0.14	1.1	$n-2$
2	1	1.3	-0.3	0.09		$= \frac{1.1}{3}$
3	2	2	0	0		
4	2	2.7	-0.7	0.49		$= 0.364$
5	4	3.4	0.4	0.36		

$H_0: \hat{\beta}_1 = 0$ vs $H_1: \hat{\beta}_1 \neq 0$ two sided t-test

$t_c = t_{\alpha/2, n-2} = t_{0.025, 3} = 3.18$, we reject H_0 if $|t| > |t_c|$

$$t = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{MSRES/S_{xx}}} = \frac{0.7}{\sqrt{\frac{0.364}{10}}} = 3.654$$

Hence, we reject H_0 and thus there is a significant relationship at 5% level.

$$\text{For } \hat{\beta}_0: t = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{MSRES} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}} = \frac{-0.1}{\sqrt{0.364} \sqrt{\frac{1}{5} + \frac{3^2}{10}}} = -0.1574$$

* To get the p-value: $P(|t|) = 2P(|t| > 3.18)$

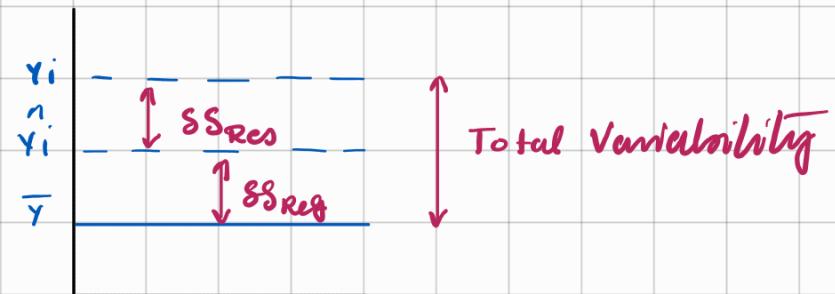
- 1) Look for $df = n-2$
- 2) Check the closest two values to the calculated t in row $n-2$
- 3) Take the average of their respective α
- 4) Multiply by 2.

$$\text{For } \hat{\beta}_1: 2 \cdot \left(\frac{0.02 + 0.015}{2} \right) = 0.035 < 0.05, \text{ we reject } H_0$$

$$\text{For } \hat{\beta}_0: 2 \cdot \left(\frac{0.48 + 0.40}{2} \right) = 0.88 > 0.05, \text{ cannot reject } H_0$$

$$\text{standard error} = \frac{\text{estimate}}{P\text{ value}} = \frac{0.7}{0.085} = 20$$

c) Build an Anova table to test $H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$



Total variability at the data = $\sum_{i=1}^n (y_i - \bar{y})^2$

$$\begin{aligned} SS_T &= \sum (y_i - \bar{y})^2 \\ &= \sum (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 \\ &= \sum (y_i - \hat{y}_i)^2 + \sum (\hat{y}_i - \bar{y})^2 + 2 \sum (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) \\ &= SS_{\text{Res}} + SS_{\text{Reg}} \end{aligned}$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i = \bar{y} + \hat{\beta}_1 (x_i - \bar{x})$$

$$\begin{aligned} \sum_{i=1}^n (y_i - \bar{y})(y_i - \hat{y}_i) &= \sum_{i=1}^n \hat{\beta}_1 (x_i - \bar{x}) ((y_i - \bar{y}) - \hat{\beta}_1 (x_i - \bar{x})) \\ &= \hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) - \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \hat{\beta}_1 S_{xy} - \hat{\beta}_1^2 S_{xx} \\ &= \hat{\beta}_1 S_{xy} - \hat{\beta}_1 S_{xy}, \quad \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \\ &= 0 \end{aligned}$$

SS_{Res} : Deviations of the predicted value from the actual value.

SS_{Reg} : Deviations of the predicted value to the mean

$$\begin{aligned}
 SS_{\text{Reg}} &= \sum (\hat{y}_i - \bar{Y})^2 \\
 &= \sum (\hat{\beta}_0 + \hat{\beta}_1 x_i - \hat{\beta}_0 - \hat{\beta}_1 \bar{x})^2 \\
 &= \sum \hat{\beta}_1^2 (x_i - \bar{x})^2 \\
 &= \hat{\beta}_1^2 \sum (x_i - \bar{x})^2 \\
 &= \hat{\beta}_1^2 S_{xx} \\
 &= \hat{\beta}_1^2 S_{xy}
 \end{aligned}$$

$$SS_{\text{Res}} = \sum (y_i - \hat{y}_i)^2 = \sum e_i^2 \sim \chi^2_{n-2} \rightarrow \text{normal equations}$$

$SS_T : \sum (y_i - \bar{Y}) = 0$ is the only restriction : $n-1$ DF
 So, SS_{Reg} has only one degree of freedom $\sim \chi^2$,

$$\mathbb{E}(MS_{\text{Res}}) = \sigma^2$$

$$\begin{aligned}
 \mathbb{E}(MS_{\text{Reg}}) &= \mathbb{E}(\hat{\beta}_1^2 S_{xx}) = S_{xx} (\text{Var}(\hat{\beta}_1) + \mathbb{E}(\hat{\beta}_1)^2) \\
 \mathbb{E}(\hat{\beta}_1^2 S_{xy}) &= S_{xx} \left(\frac{\sigma^2}{S_{xx}} + \hat{\beta}_1^2 \right) \\
 &= \sigma^2 + S_{xx} \hat{\beta}_1^2
 \end{aligned}$$

ANOVA TABLE :

Source of Variability	DF	SS	MS	F-Value
Regression	1	4.9	4.9	$F = \frac{MS_{\text{Reg}}}{MS_{\text{Res}}} \sim F_{1,n-2}$
Residual	3	1.1	0.366	
Total	4	6		$F = \frac{4.9}{0.366} = 13.38$

$F > F_{1, n-2, 0.05} = 10.13$ so we reject H_0 and thus f_1 is statistically significant and the model is significant

$$SST = \sum_{i=1}^5 (Y_i - \bar{Y})^2 = S_{YY} = 6$$

$$SS_{RES} = \sum_{i=1}^5 (Y_i - \hat{Y}_i)^2 = 0.16 + 0.09 + 0.49 + 0.36 = 1.1$$

$$\Rightarrow SS_{Reg} = SST - SS_{RES} = 4.9$$

Alternatively, $SS_{Reg} = \sum (Y_i - \hat{Y}_i)^2$

$$= (0.6 - 2)^2 + (1.3 - 2)^2 + (2 - 2)^2 + (2.7 - 2)^2 + (3.4 - 2)^2 = 4.9$$

$$MS_{Reg} = \frac{SS_{Reg}}{1} = \frac{4.9}{1} = 4.9$$

$$MS_{Res} = \frac{SS_{Res}}{3} = \frac{1.1}{3} = 0.366$$

d) Check the accuracy of the regression model

Measure of "Goodness of Fit":

Coefficient of determination: R^2

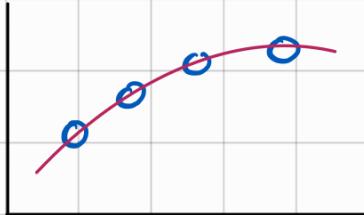
R^2 measures the proportion of variability in response variable that is explained by the regression model.

$$R^2 = \frac{SS_{Reg}}{SST} = 1 - \frac{SS_{Res}}{SST}, \text{ maximum } SS_{Reg}$$

minimum SS_{Res}

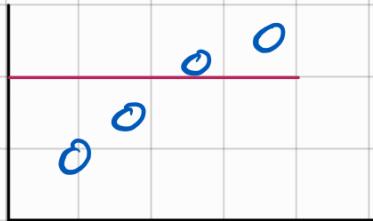
$$r_{xy} = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var} x \text{Var} y}}$$

- $R^2 = 1 : SS_{\text{Res}} = SS_{\text{Reg}}$ Perfect fit



$SS_{\text{Res}} = 0$ so all the variability is explained by the model (var of the data)

- $R^2 = 0 : SS_{\text{Res}} = SS_{\text{T}}$ $\Rightarrow \sum (y_i - \hat{y})^2 = \sum (y_i - \bar{y})^2$
 $\Rightarrow \hat{y}_i = \bar{y}$



No relationship so $\hat{\beta}_1 = 0$ then $S_{xy} = 0$
 $R^2 = \frac{SS_{\text{Res}}}{SS_{\text{T}}} = \frac{\hat{\beta}_1^2 S_{xx} + 0}{SS_{\text{T}}} = 0 \Rightarrow \hat{\beta}_1 = 0$

$$R^2 = \frac{SS_{\text{Reg}}}{SS_{\text{T}}} = \frac{4.9}{6} = 0.816 = 0.82$$

Therefore, 82% of the total variability of the sales amount is explained by advertisement. It is not 100% because, it would an over fitted model and there might be some other variables that are useful to explain the remaining variability of the model but was not given e.g # of outlets & # of workers.

e) Find confidence interval for β_1 .

$$\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{S_{xx}}), \text{ Normal as } Y \sim N(\beta_0 + \beta_1 x, \sigma^2)$$

$$\text{By CLT: } z = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\sigma^2 / s_{xx}}} \sim N(0, 1)$$

$$t = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{MSRes}{s_{xx}}}} \sim t_{\alpha/2, n-2}, \quad \sigma^2 = E(MSRes)$$

$$\text{Thus, } P\left(\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{MSRes}{s_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{MSRes}{s_{xx}}}\right) = 1 - \alpha$$

$$P\left(0.7 - t_{0.025, 3} \sqrt{\frac{0.364}{10}} \leq \beta_1 \leq 0.7 + t_{0.025, 3} \sqrt{\frac{0.364}{10}}\right) = 0.95$$

$$\beta_1 \in [0.091, 1.31], \quad t_{0.025, 3} = 3.182$$

f) Estimate the mean sales amount when advertising cost is 4 (40K \$) at 0.05 level.

$$\hat{Y} = -0.1 + 0.7X, \quad X = 4 \Rightarrow \hat{Y} = 2.7$$

Interval estimation of mean response given $X = m_0$:

$$E(Y | X = m_0) = \hat{\beta}_0 + \hat{\beta}_1 m_0 = \hat{Y}$$

$$\begin{aligned} \text{Var}(\hat{\beta}_0 + \hat{\beta}_1 m_0) &= \text{Var}(\bar{Y} + \hat{\beta}_1 (m_0 - \bar{X})) \\ &= \text{Var} \bar{Y} + (m_0 - \bar{X})^2 \text{Var} \hat{\beta}_1 - 2 \text{Cov}(\bar{Y}, \hat{\beta}_1 (m_0 - \bar{X})) \\ &= \frac{\sigma^2}{n} + (m_0 - \bar{X})^2 \cdot \frac{\sigma^2}{s_{xx}} - 2(m_0 - \bar{X}) \text{Cov}(\bar{Y}, \hat{\beta}_0) \\ &= \sigma^2 \left(\frac{1}{n} + \frac{(m_0 - \bar{X})^2}{s_{xx}} \right) \end{aligned}$$

independent
 \bar{Y} and $\hat{\beta}_0$

$$E(Y | X = m_0) \sim N(\hat{Y}, \hat{\sigma}^2 / (1 + (m_0 - \bar{X})^2))$$

$$\mathbb{E}(Y|X=n_0) \sim N(\beta_0 + \beta_1 n_0, s^2 = \left(\frac{1}{n} + \frac{(n_0 - \bar{x})^2}{s_{xx}} \right))$$

$$\mathbb{E}(Y|X=n_0) \in [\hat{\beta}_0 + \hat{\beta}_1 n_0 - t_{\frac{\alpha}{2}, n-2} \sqrt{MSRes \left(\frac{1}{n} + \frac{(n_0 - \bar{x})^2}{s_{xx}} \right)}, \hat{\beta}_0 + \hat{\beta}_1 n_0 + t_{\frac{\alpha}{2}, n-2} \sqrt{MSRes \left(\frac{1}{n} + \frac{(n_0 - \bar{x})^2}{s_{xx}} \right)}]$$

$$\mathbb{E}(Y|X=n_0) \in [2.7 - 3.812 \sqrt{0.366 \left(\frac{1}{5} + \frac{(4-3)^2}{10} \right)}, 2.7 + 3.812 \sqrt{0.366 \left(\frac{1}{5} + \frac{(4-3)^2}{10} \right)}]$$

$$\mathbb{E}(Y|X=n_0) \in [1.65, 3.75]$$

$$P(1.65 \leq \mathbb{E}(Y|X=n_0) \leq 3.75) = 0.95$$